

Lecture 20: Algorithms and Implications of Calibration

November 6, 2025

*Lecturer: Nika Haghtalab**Readings: Foster and Hart [2021]**Scribe: Gongyu Chen, William Yan*

1 Lecture Overview

Last time, we introduced the notion of calibrated forecasting. For a sequence of calibrated predictions $c_1, \dots, c_T \in [0, 1]$ that predicts the *outcomes* $a_1, \dots, a_T \in \{0, 1\}$, we defined the *expected* calibration error by:

$$\text{CalErr} := \sum_{p \in [0,1]} \frac{n_T(p)}{T} |\bar{a}_T(p) - p|,$$

where

$$n_T(p) = \sum_{t=1}^T \mathbb{1}(c_t = p), \quad \bar{a}_T(p) = \frac{1}{n_T(p)} \sum_{t=1}^T \mathbb{1}(c_t = p) \cdot a_t$$

We then discussed that a forecaster can indeed be calibrated even if the outcomes are adversarial decided, as characterized by the following two theorems. Informally, we derived:

- Theorem A: There exists a randomized algorithm for forecasting c_1, \dots, c_T that satisfies for any $\varepsilon > 0$:

$$\lim_{T \rightarrow \infty} \text{CalErr} < \varepsilon.$$

- Theorem B: Given that the adversarial's mixed strategy over a_1, \dots, a_T , there exists c_1, \dots, c_T such that $\lim_{T \rightarrow \infty} \text{CalErr} < \varepsilon$.

To prove Theorem A, we used the min-max theorem (that holds when the mixed strategies are finite) and showed that Theorem B implies Theorem A. Then to show Theorem B, our strategy was to define c_t by rounding the distribution of adversary's strategy $p_t = \mathbb{E}[a_t \mid \text{history}_{t-1}]$ to the discretized grid $D = \{\frac{1}{2N}, \frac{3}{2N}, \dots, \frac{2N-1}{2N}\}$.

The problem with the above min-max proof is that it does not give an *implementable* algorithm for Theorem A: The algorithm in Theorem B depends on unseen adversarial's distribution over a_t . In today's lecture, we are going to prove Theorem A in a constructive way. There exist several algorithms that achieve the claim of Theorem A, and we will see a simple one by Foster and Hart [2021].

2 Constructive Proof of Theorem A

We recall the following setup. At each time $t = 1, \dots, T$:

- An outcome $a_t \in \{0, 1\}$ is realized.
- Before observing a_t , the forecaster issues a (discretized) probability forecast

$$c_t \in D, \quad D := \left\{0, \frac{1}{N}, \frac{2}{N}, \dots, \frac{N}{N}\right\},$$

for some fixed integer $N \geq 1$. Notice that this is a different discretization from the last time.

For any forecast level $d \in D$ and time $t \geq 1$, define the number of times the forecaster predicts d up to time $t - 1$ by

$$n_{t-1}(d) := \sum_{\tau=1}^{t-1} \mathbb{1}(c_\tau = d), \quad (1)$$

and define the corresponding empirical success frequency

$$\bar{a}_{t-1}(d) := \frac{1}{n_{t-1}(d)} \sum_{\tau=1}^{t-1} \mathbb{1}(c_\tau = d) a_\tau. \quad (2)$$

For $d \in D$ and $t \geq 1$, the *normalized gap* is the unsigned error:

$$e_{t-1}(d) := \bar{a}_{t-1}(d) - d. \quad (3)$$

which, recall the notations from the last time, equals $G_{t-1}(d)/n_{t-1}(d)$ for the unsigned gap $G(d)$ normalized by $n(d)$. Finally, the calibration error is

$$K_T := \frac{1}{T} \sum_{d \in D} n_T(d) |e_T(d)|. \quad (4)$$

2.1 Algorithm

At time $t = 1, \dots, T$, given sofar observations $(a_1, c_1), \dots, (a_{t-1}, c_{t-1})$. compute $e_{t-1}(d)$ for all $d \in D$. Then there are two cases on $e_{t-1}(d)$.

- **Case 1:** If there exists $d \in D$ such that $e_{t-1}(d) = 0$, pick any such d and set

$$c_t := d.$$

For later use, we also let $y := c_t$.

- **Case 2:** Otherwise, if all $e_{t-1}(d_k) \neq 0$ for all $d \in D$. then there must exists two consecutive grid values $j-1, j \in [N]$ such that

$$e_{t-1}\left(\frac{j-1}{N}\right) > 0 > e_{t-1}\left(\frac{j}{N}\right).$$

Then let

$$c_t = \begin{cases} \frac{j-1}{N} =: y_1 & \text{with probability } p_1, \\ \frac{j}{N} =: y_2 & \text{with probability } p_2, \end{cases}$$

where

$$p_1 := \frac{|e_{t-1}(y_2)|}{|e_{t-1}(y_1)| + |e_{t-1}(y_2)|}, \quad p_2 := \frac{|e_{t-1}(y_1)|}{|e_{t-1}(y_1)| + |e_{t-1}(y_2)|}.$$

The justification for the existence of such adjacent pair of $j-1$ and j is as follows. For $d = 0$, the prediction can only underestimate the outcome, thus $e_{t-1}(0) = \bar{a}_{t-1}(0) > 0$. Similarly, for $d = 1$, the forecaster is asserting that the outcome event is happening with 100% probability, thus the prediction can only overestimate the true outcome, so $e_{t-1}(1) < 0$. The existence of such $(j-1, j)$ is thus obtained by the intermediate value theorem applied to the continuous function $e_{t-1}(d)$ with domain $[0, 1]$. Intuitively, the algorithm chooses the prediction that interpolates between y_1 and y_2 in expectation, so approximately there is no overestimation or underestimation bias.

2.2 Constructive Guarantee

Our goal is to show that the above constructive algorithm ensures the following performance guarantee.

Theorem 2.1. *For $T \geq N^3$, the algorithm in Section 2.1 guarantees that:*

$$\mathbb{E}[K_T] \leq O\left(\frac{1}{N}\right),$$

hence with $N \asymp \frac{1}{\varepsilon}$, we get $\lim_{T \rightarrow \infty} \mathbb{E}[K_T] \leq \varepsilon$ and prove the Theorem A.

The above theorem is implied by a major lemma that unpacks why the theorem works. Intuitively, a forecast c_t is calibrated if the error $e_{t-1}(c_t)$ is close to zero, with slight over- or under-estimations. To control the error, the Major Lemma 2.2 presents a condition (5) that approximately says the error of a forecast c_t is small as measured by the correlation of its past error $e_{t-1}(c_t)$ and its future error $a - c_t$ for unseen outcome $a \in \{0, 1\}$, and this correlation measure is at most the standard deviation of the past error $e_{t-1}(c_t)$.

Lemma 2.2 (Major lemma). *Given a randomized forecaster algorithm with forecasts from a finite set D , suppose that for every t , any history h_{t-1} , and both outcomes $a \in \{0, 1\}$,*

$$\mathbb{E}[e_{t-1}(c_t)(a - c_t) \mid h_{t-1}] \leq \varepsilon \cdot \mathbb{E}[|e_{t-1}(c_t)| \mid h_{t-1}], \quad (5)$$

then the calibration error satisfies:

$$\mathbb{E}[K_T] \leq \varepsilon + \tilde{O}\left(\sqrt{\frac{N}{T}}\right) \implies \lim_{T \rightarrow \infty} \mathbb{E}[K_T] \leq \varepsilon.$$

The proof of the above lemma is in the optional reading [Foster and Hart, 2021]. We are focusing on establishing why the algorithm satisfies the precondition (5). We first want to discuss the similar structure of the condition to the conditions appeared in the proof of Theorem B. Recall that in Theorem B, given the adversary's mixed strategy p_t , the forecast c_t rounds it to the grid D so $p_t - c_t \approx 1/(2N)$ is the rounding error that is very negligible for large N , and the error against the true outcome a_t satisfies:

$$\mathbb{E}[(a_t - c_t)|h_t] = 0, \quad \mathbb{E}[(a_t - c_t)^2|h_t] \approx \mathbb{E}[(a_t - p_t)^2] \leq \frac{1}{4}, \text{ since } a_t \sim \text{Bernoulli}(p_t).$$

Notice that the variance condition here is similar to the condition (5), except that we now have no control over a_t , so we expect the condition (5) to hold for arbitrary outcome $a \in \{0, 1\}$. Further, the rounding bias $p_t - c_t$ here is also similar to the normalized gap $e_t(c_t)$ defined for Lemma 2.2.

We now formally prove that the Algorithm in Section 2.1 satisfies the condition (5).

Proof. • Claim 1:

$$\mathbb{E}[e_{t-1}(c_t)|h_{t-1}] = 0$$

It is easy to see that if we are in Case 1 of the Algorithm, $e_{t-1}(c_t) = 0$ for arbitrary $c_t \in D$.

For Case 2, we can expand:

$$\mathbb{E}[e_{t-1}(c_t) | h_{t-1}] = p_1 \cdot e_{t-1}(y_1) + p_2 \cdot e_{t-1}(y_2) = 0$$

where we use the definition of p_1, p_2 , and the fact that $e_{t-1}(y_1) > 0 > e_{t-1}(y_2)$.

• Claim 2: For all $a \in \{0, 1\}$,

$$\mathbb{E}[e_{t-1}(c_t) (a - c_t) | h_{t-1}] \leq \frac{1}{2N} \mathbb{E}[|e_{t-1}(c_t)| | h_{t-1}].$$

Define the *anchor*:

$$\bar{y} = \begin{cases} y & \text{Case 1} \\ \frac{1}{2}(y_1 + y_2) & \text{Case 2} \end{cases}$$

then by inserting the anchor, we can consider the following decomposition:

$$\mathbb{E}[e_{t-1}(c_t) \underbrace{(a - \bar{y})}_{(1)} + \underbrace{\bar{y} - c_t}_{(2)} | h_{t-1}]$$

where

$$(1) = \mathbb{E}[e_{t-1}(c_t) (a - \bar{y}) \mid h_{t-1}] = (a - \bar{y}) \cdot \mathbb{E}[e_{t-1}(c_t) \mid h_{t-1}] = 0, \quad (6)$$

from the facts that a is fixed and \bar{y} has no additional randomness conditioned on the history h_{t-1} ; and

$$(2) = \mathbb{E}[e_{t-1}(c_t) (\bar{y} - c_t) \mid h_{t-1}] \leq |\bar{y} - c_t| \cdot \mathbb{E}[|e_{t-1}(c_t)|] \leq \frac{1}{2N} \cdot \mathbb{E}[|e_{t-1}(c_t)|]$$

from the discretization. Combining the two terms recovers the Claim 2.

(*Proof of Theorem 2.1*) Now if $T = N^3$, by the Major Lemma and Claim 2, we have:

$$\lim_{T \rightarrow \infty} \mathbb{E}[K_T] \leq \frac{1}{2N}.$$

□

2.3 Generalization to multiclass calibration

Fix an outcome alphabet $[K] := \{1, 2, \dots, K\}$. We observe a sequence of outcomes

$$a_1, a_2, \dots, a_T \in [K].$$

At each time t , a forecaster outputs a vector from the probability simplex

$$c_t \in \Delta_K := \left\{ p \in \mathbb{R}_{\geq 0}^K : \sum_{k=1}^K p_k = 1 \right\},$$

so c_t is a probability distribution over K items.

Analogue to the binary case, we say that c_1, \dots, c_T is ε -calibrated w.r.t. a_1, \dots, a_T if

$$\lim_{T \rightarrow \infty} \sum_{p \in \Delta_K} \frac{n_T(p)}{T} \|\bar{a}_T(p) - p\|_1 \leq \varepsilon.$$

Here, the empirical quantities are again:

$$n_T(p) := \sum_{t=1}^T \mathbb{1}\{c_t = p\},$$

$$\bar{a}_T(p) = \frac{1}{n_T(p)} \sum_{t=1}^T \mathbb{1}\{c_t = p\} \cdot \delta(a_t)$$

where $\delta(a_t)$ is the standard basis vector with a 1 at coordinate a_t and 0 elsewhere.

We compare the guarantees for multiclass and binary calibrations:

Setting	How long to run (T)	Remarks / reference
1. Binary	$T \gtrsim \varepsilon^{-3}$	Theorem A/B from last class
2. K -class	$T \gtrsim \left(\frac{1}{\varepsilon^3}\right)^{\Theta(K)}$	Generalization of min-max (Thm B) with grid $D_K = \underbrace{D \times D \cdots \times D}_{K \text{ times}}$
3. K -class	$T \gtrsim (K)^{O(1/\varepsilon^2)}$	[Peng, 2025, Fishelson et al., 2025]

Table 1: Summary of calibration-error bounds (“how long to run”)

3 From calibration to no swap regret

One reason that we care about calibrated forecasting is that decisions that are made based on such forecasts are of high quality. In particular, that they incur no-swap regret.

The high-level reason that the calibration is related to the swap regret is as follows: If a *calibrated* forecaster predicts an event’s likelihood is 70%, then we can trust the prediction and assume the event is truly going to happen with probability 70%; on the other hand, for a forecaster that is not calibrated and consistently overestimates the outcomes, if they predicts 70%, we should *swap* to a downgraded prediction (say 60%). This effectively means that following the advice of a calibrated predictor leads to *no swap regret*. We formalize this below.

Best response to calibrated forecasts. Now let an agent have a finite action set S , and suppose on each round t :

- Nature (or other players) produces outcome $a_t \in [K]$,
- The agent *observes the forecast* $c_t \in \mathcal{P}$,
- The agent chooses an action $s_t \in S$.

The agent’s utility is given by

$$u : S \times [K] \rightarrow [0, 1], \quad u(s, a) \in [0, 1].$$

Given ε -calibrated forecasts c_1, \dots, c_T of unseen outcomes $a_1, \dots, a_T \in [K]$, the agent’s best response is

$$s_t \in \text{BR}(c_t) := \arg \max_{s \in S} \mathbb{E}_{a \sim c_t} [u(s, a)] \quad (7)$$

Theorem 3.1 (calibration \Rightarrow no swap regret). *If c_1, \dots, c_T is ε -calibrated for a_1, \dots, a_T , then the sequence of actions s_1, \dots, s_T played by the agent has swap regrets at most εT , i.e.*

$$\max_{\varphi: S \rightarrow S} \sum_{t=1}^T u(\varphi(s_t), a_t) - \sum_{t=1}^T u(s_t, a_t) \leq \varepsilon T.$$

Proof. Fix any swap function $\varphi : S \rightarrow S$. Define the swap regret w.r.t. φ :

$$R_\varphi^{\text{swap}} := \sum_{t=1}^T u(\varphi(s_t), a_t) - \sum_{t=1}^T u(s_t, a_t).$$

Because $s_t = \text{BR}(c_t)$, we can rewrite this as

$$\begin{aligned} R_\varphi^{\text{swap}} &= \sum_{t=1}^T u(\varphi(\text{BR}(c_t)), a_t) - \sum_{t=1}^T u(\text{BR}(c_t), a_t) \\ &= \sum_{p \in \Delta_K} \sum_{t=1}^T \mathbb{1}(c_t = p) (u(\varphi(\text{BR}(p))), a_t) - u(\text{BR}(p), a_t). \end{aligned}$$

Introduce the shorthand

$$\Delta(p) := \mathbb{E}_{a \sim \bar{a}_T(p)} [u(\varphi(\text{BR}(p)), a) - u(\text{BR}(p), a)],$$

i.e. the expected gain from swapping $\text{BR}(p)$ to $\varphi(\text{BR}(p))$ when the outcome distribution is the empirical one corresponding to forecast p . Then

$$R_\varphi^{\text{swap}} = \sum_{p \in \Delta_K} n_T(p) \Delta(p).$$

So it suffices to bound $\Delta(p)$ by the calibration error of bin p .

Claim. For every $p \in \Delta_K$,

$$\Delta(p) \leq \|\bar{a}_T(p) - p\|_1.$$

Proof of claim. We use the following basic inequality: if $f : [K] \rightarrow [0, 1]$ and μ, ν are distributions on $[K]$, then

$$|\mathbb{E}_{a \sim \mu}[f(a)] - \mathbb{E}_{a \sim \nu}[f(a)]| \leq \|\mu - \nu\|_1. \quad (8)$$

1. Apply (8) to $f(a) = u(\varphi(\text{BR}(p)), a)$, $\mu = \bar{a}_T(p)$, $\nu = p$:

$$\left| \mathbb{E}_{a \sim \bar{a}_T(p)} u(\varphi(\text{BR}(p)), a) - \mathbb{E}_{a \sim p} u(\varphi(\text{BR}(p)), a) \right| \leq \|\bar{a}_T(p) - p\|_1.$$

2. Apply (8) to $f(a) = u(\text{BR}(p), a)$, $\mu = \bar{a}_T(p)$, $\nu = p$:

$$\left| \mathbb{E}_{a \sim \bar{a}_T(p)} u(\text{BR}(p), a) - \mathbb{E}_{a \sim p} u(\text{BR}(p), a) \right| \leq \|\bar{a}_T(p) - p\|_1.$$

3. By definition of best response,

$$\mathbb{E}_{a \sim p} u(\text{BR}(p), a) \geq \mathbb{E}_{a \sim p} u(\varphi(\text{BR}(p)), a).$$

Now write $\Delta(p)$ and add/subtract the expectations under $a \sim p$:

$$\begin{aligned}\Delta(p) &= \mathbb{E}_{a \sim \bar{a}_T(p)} u(\varphi(\text{BR}(p)), a) - \mathbb{E}_{a \sim \bar{a}_T(p)} u(\text{BR}(p), a) \\ &= [\mathbb{E}_{a \sim \bar{a}_T(p)} u(\varphi(\text{BR}(p)), a) - \mathbb{E}_{a \sim p} u(\varphi(\text{BR}(p)), a)] \\ &\quad + [\mathbb{E}_{a \sim p} u(\varphi(\text{BR}(p)), a) - \mathbb{E}_{a \sim p} u(\text{BR}(p), a)] \\ &\quad + [\mathbb{E}_{a \sim p} u(\text{BR}(p), a) - \mathbb{E}_{a \sim \bar{a}_T(p)} u(\text{BR}(p), a)].\end{aligned}$$

The middle bracket is ≤ 0 by the best-response property. The first and third brackets are each bounded in absolute value by $\|\bar{a}_T(p) - p\|_1$ by the two applications of (8) above. Hence

$$\Delta(p) \leq \|\bar{a}_T(p) - p\|_1 + 0 + \|\bar{a}_T(p) - p\|_1 \leq 2\|\bar{a}_T(p) - p\|_1.$$

□

Returning to the main proof, we have

$$R_\varphi^{\text{swap}} = \sum_{p \in \Delta_K} n_T(p) \Delta(p) \leq \sum_{p \in \Delta_K} n_T(p) \|\bar{a}_T(p) - p\|_1 \leq \varepsilon T,$$

by the assumed ε -calibration. Since this holds for every $\varphi : S \rightarrow S$, the max over φ is also $\leq \varepsilon T$. □

Crucially, the proof shows that if the agent is best-responding to the forecasts calibrated by the empirical distribution of outcomes, they won't want to swap out of the best response.

References

- Maxwell Fishelson, Noah Golowich, Mehryar Mohri, and Jon Schneider. High-dimensional calibration from swap regret. *arXiv preprint arXiv:2505.21460*, 2025.
- Dean P Foster and Sergiu Hart. Forecast hedging and calibration. *Journal of Political Economy*, 129(12):3447–3490, 2021.
- Binghui Peng. High dimensional online calibration in polynomial time. *arXiv preprint arXiv:2504.09096*, 2025.